# CS252 Graduate Computer Architecture Spring 2014
# Lecture 8: Advanced Out-of-Order Superscalar Designs Part-II

Krste Asanovic

**krste@eecs.berkeley.edu**
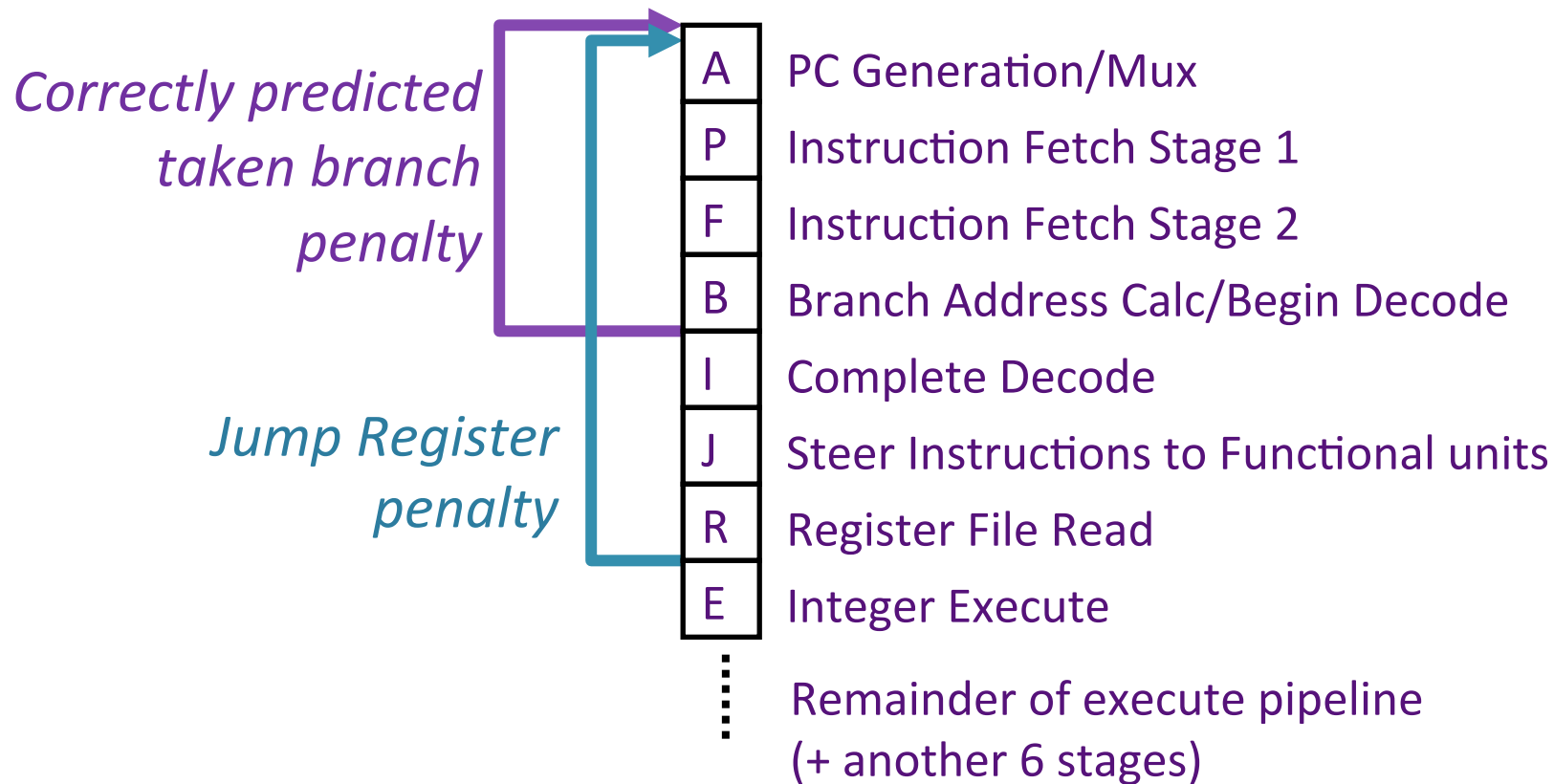
**http://inst.eecs.berkeley.edu/~cs252/fa15**

# Last Time in Lecture 7

- Unified Physical Register Design for OoO superscalar
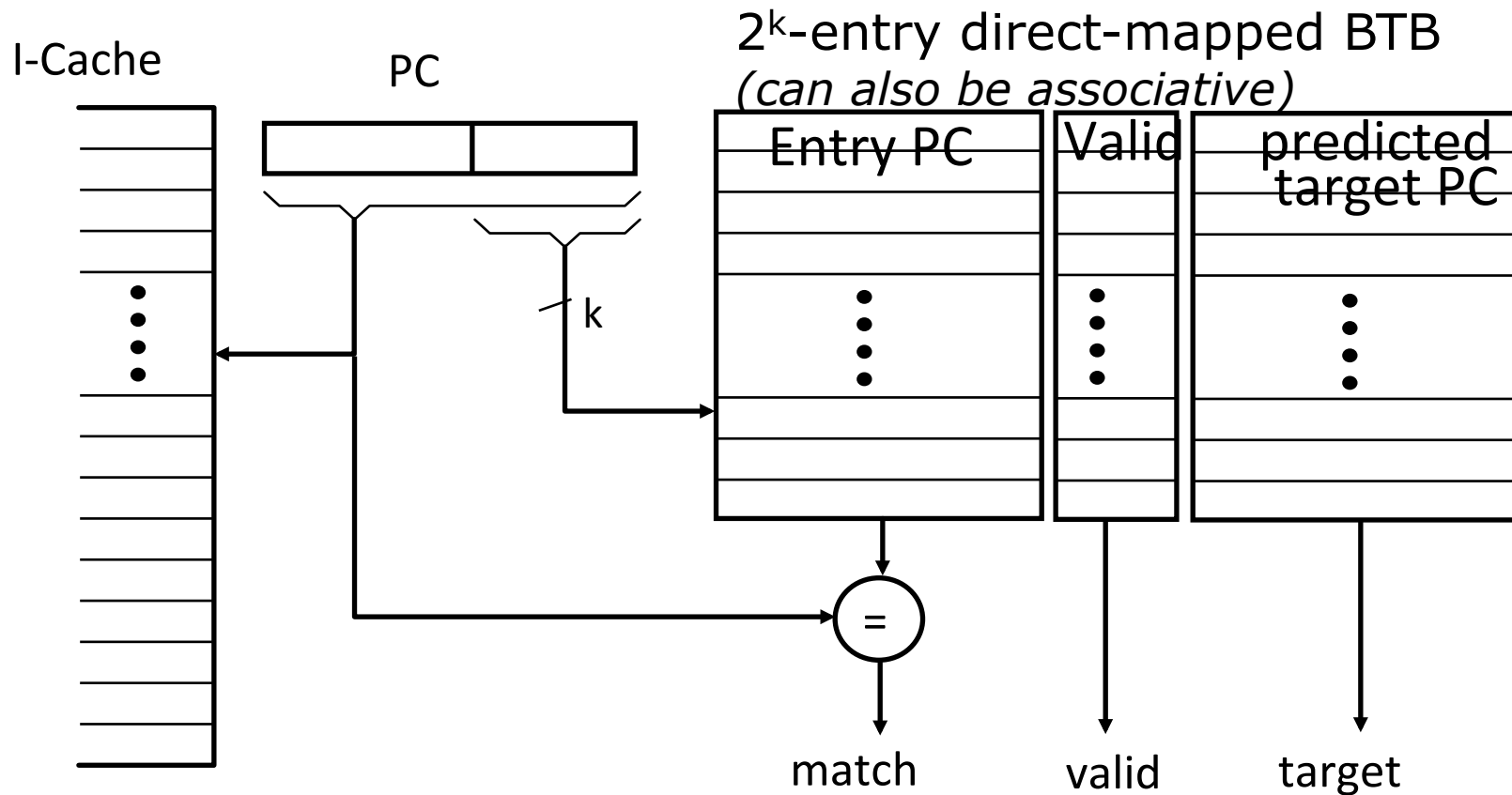- Branch History Table Branch Predictors

# Limitations of BHTs

Only predicts branch direction. Therefore, cannot redirect fetch stream until after branch target is determined.

*Correctly predicted taken branch penalty*

*Jump Register penalty*

| | |
|---|---|
| A | PC Generation/Mux |
| P | Instruction Fetch Stage 1 |
| F | Instruction Fetch Stage 2 |
| B | Branch Address Calc/Begin Decode |
| I | Complete Decode |
| J | Steer Instructions to Functional units |
| R | Register File Read |
| E | Integer Execute |

Remainder of execute pipeline
(+ another 6 stages)

*UltraSPARC-III fetch pipeline*

# Branch Target Buffer (BTB)

I-Cache      PC      $2^k$-entry direct-mapped BTB
*(can also be associative)*

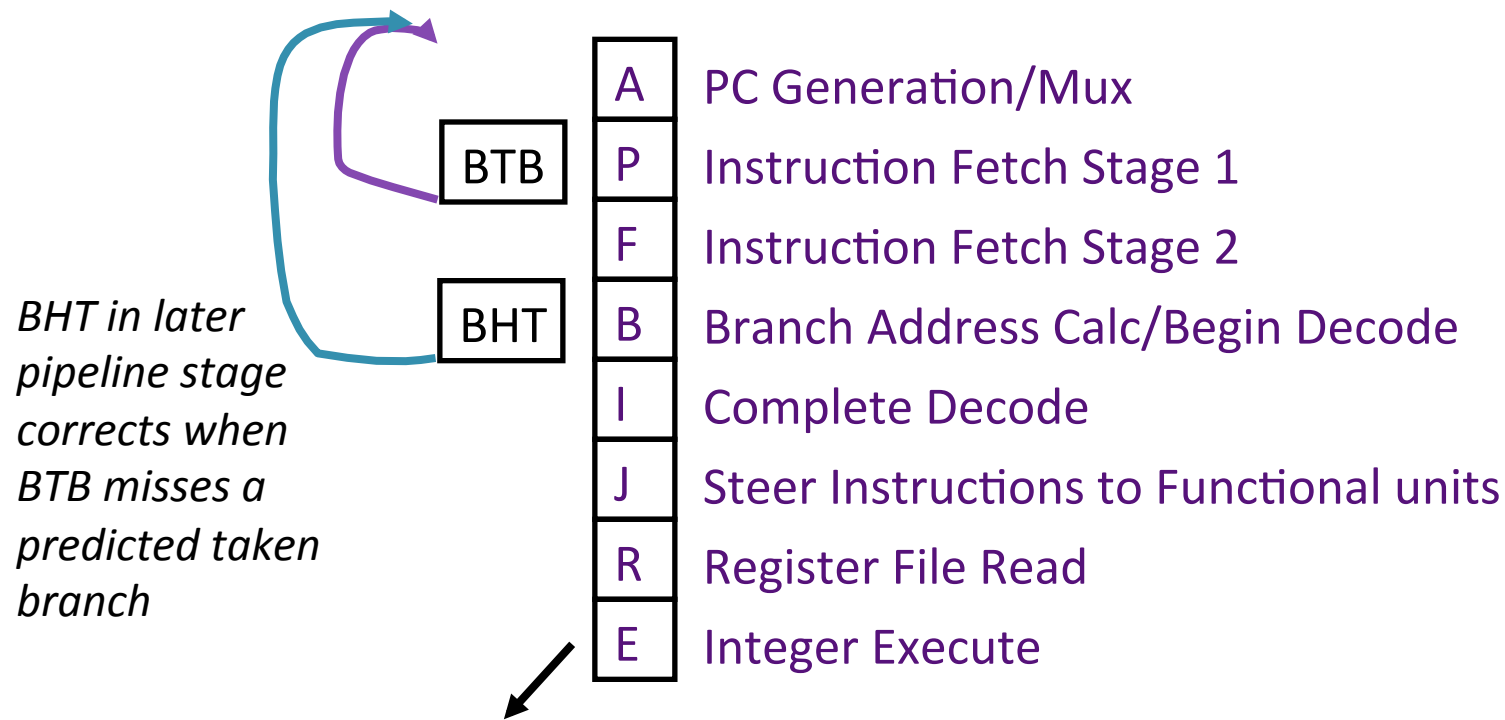| Entry PC | Valid | predicted target PC |
|---|---|---|

k

match     valid     target

- Keep both the branch PC and target PC in the BTB
- PC+4 is fetched if match fails
- Only *taken* branches and jumps held in BTB
- Next PC determined *before* branch fetched and decoded

# Combining BTB and BHT

- BTB entries are considerably more expensive than BHT, but can redirect fetches at earlier stage in pipeline and can accelerate indirect branches (JR)
- BHT can hold many more entries and is more accurate

| | Stage |
|---|---|
| A | PC Generation/Mux |
| P | Instruction Fetch Stage 1 |
| F | Instruction Fetch Stage 2 |
| B | Branch Address Calc/Begin Decode |
| I | Complete Decode |
| J | Steer Instructions to Functional units |
| R | Register File Read |
| E | Integer Execute |

BTB

BHT

*BHT in later pipeline stage corrects when BTB misses a predicted taken branch*

*BTB/BHT only updated after branch resolves in E stage*

# Uses of Jump Register (JR)

- Switch statements (jump to address of matching case)

  BTB works well if same case used repeatedly

- Dynamic function call (jump to run-time function address)

  BTB works well if same function usually called, (e.g., in C++ programming, when objects have same type in virtual function call)

- Subroutine returns (jump to return address)

  BTB works well if usually return to the same place

  ⇒ *Often one function called from many distinct call sites!*

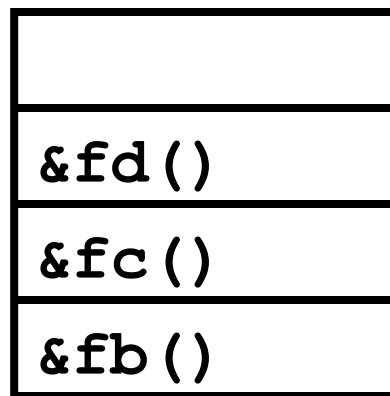  How well does BTB work for each of these cases?

# Subroutine Return Stack

Small structure to accelerate JR for subroutine returns, typically much more accurate than BTBs.

```
fa() { fb(); }
fb() { fc(); }
fc() { fd(); }
```
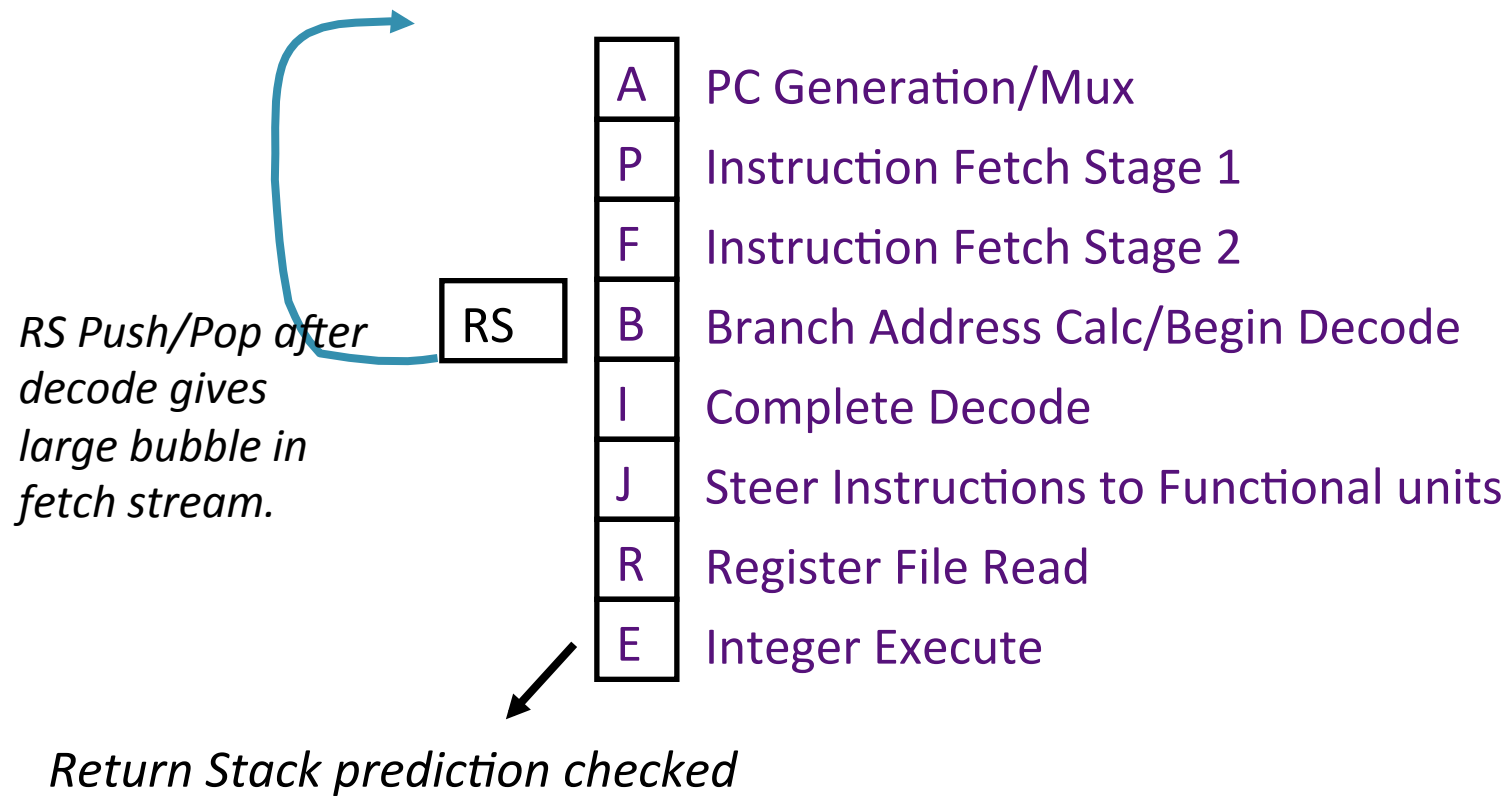
*Push call address when function call executed*

*Pop return address when subroutine return decoded*

| |
|---|
| |
| `&fd()` |
| `&fc()` |
| `&fb()` |

*k entries (typically k=8-16)*

# Return Stack in Pipeline

- How to use return stack (RS) in deep fetch pipeline?
- Only know if subroutine call/return at decode

*RS Push/Pop after decode gives large bubble in fetch stream.*

RS

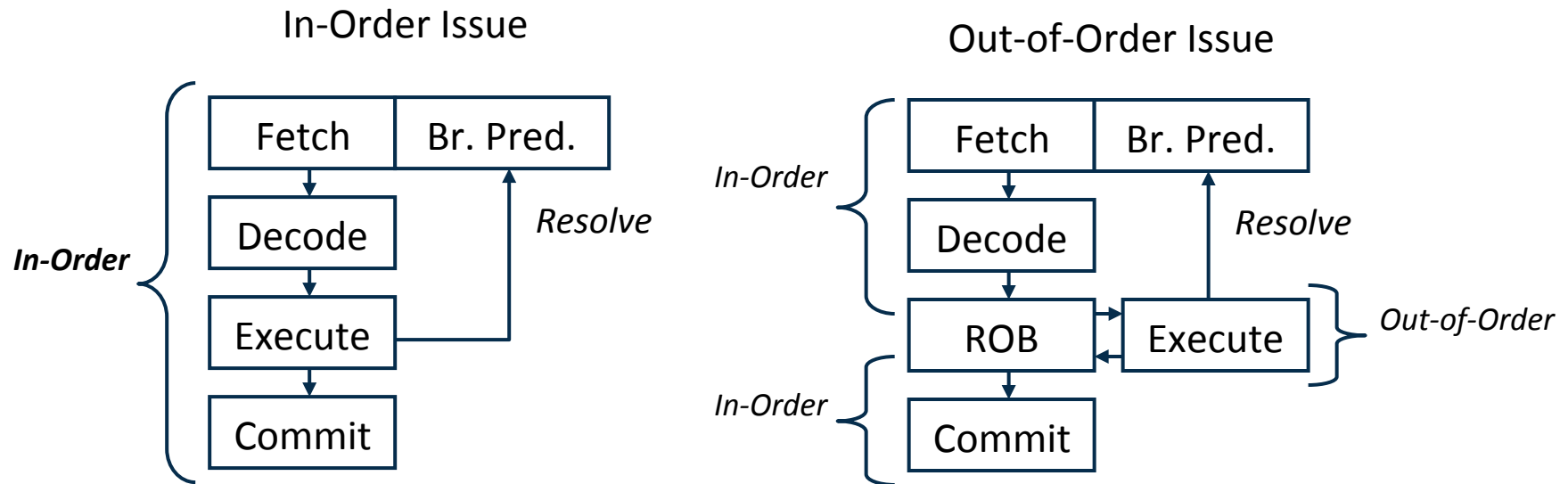| A | PC Generation/Mux |
| P | Instruction Fetch Stage 1 |
| F | Instruction Fetch Stage 2 |
| B | Branch Address Calc/Begin Decode |
| I | Complete Decode |
| J | Steer Instructions to Functional units |
| R | Register File Read |
| E | Integer Execute |

*Return Stack prediction checked*

# Return Stack in Pipeline

- Can remember whether PC is subroutine call/return using BTB-like structure
- Instead of target-PC, just store push/pop bit

| | | |
|---|---|---|
| RS | A | PC Generation/Mux |
| | P | Instruction Fetch Stage 1 |
| | F | Instruction Fetch Stage 2 |
| | B | Branch Address Calc/Begin Decode |
| | I | Complete Decode |
| | J | Steer Instructions to Functional units |
| | R | Register File Read |
| | E | Integer Execute |

*Push/Pop before instructions decoded!*

*Return Stack prediction checked*

# In-Order vs. Out-of-Order Branch Prediction

### In-Order Issue

```
In-Order {  [ Fetch | Br. Pred. ]
               |                    ↑
               ↓              Resolve
            [ Decode ]
               |
               ↓
            [ Execute ]──────────
               |
               ↓
            [ Commit ]
```

### Out-of-Order Issue

```
In-Order {  [ Fetch | Br. Pred. ]
               |                    ↑
               ↓              Resolve
            [ Decode ]
               |
               ↓
            [ ROB ] ←→ [ Execute ]  } Out-of-Order
In-Order {     |
               ↓
            [ Commit ]
```

- Speculative fetch but not speculative execution - branch resolves before later instructions complete
- Completed values held in bypass network until commit

- Speculative execution, with branches resolved after later instructions complete
- Completed values held in rename registers in ROB or unified physical register file until commit

• Both styles of machine can use same branch predictors in front-end fetch pipeline, and both can execute multiple instructions per cycle

• Common to have 10-30 pipeline stages in either style of design

# InO vs. OoO Mispredict Recovery

- In-order execution?
  - Design so no instruction issued after branch can write-back before branch resolves
  - Kill all instructions in pipeline behind mispredicted branch
- Out-of-order execution?
  - Multiple instructions following branch in program order can complete before branch resolves
  - A simple solution would be to handle like precise traps
    - Problem?

# Branch Misprediction in Pipeline



- Can have multiple unresolved branches in ROB
- Can resolve branches out-of-order by killing all the instructions in ROB that follow a mispredicted branch
- MIPS R10K uses four mask bits to tag instructions that are dependent on up to four speculative branches
- Mask bits cleared as branch resolves, and reused for next branch
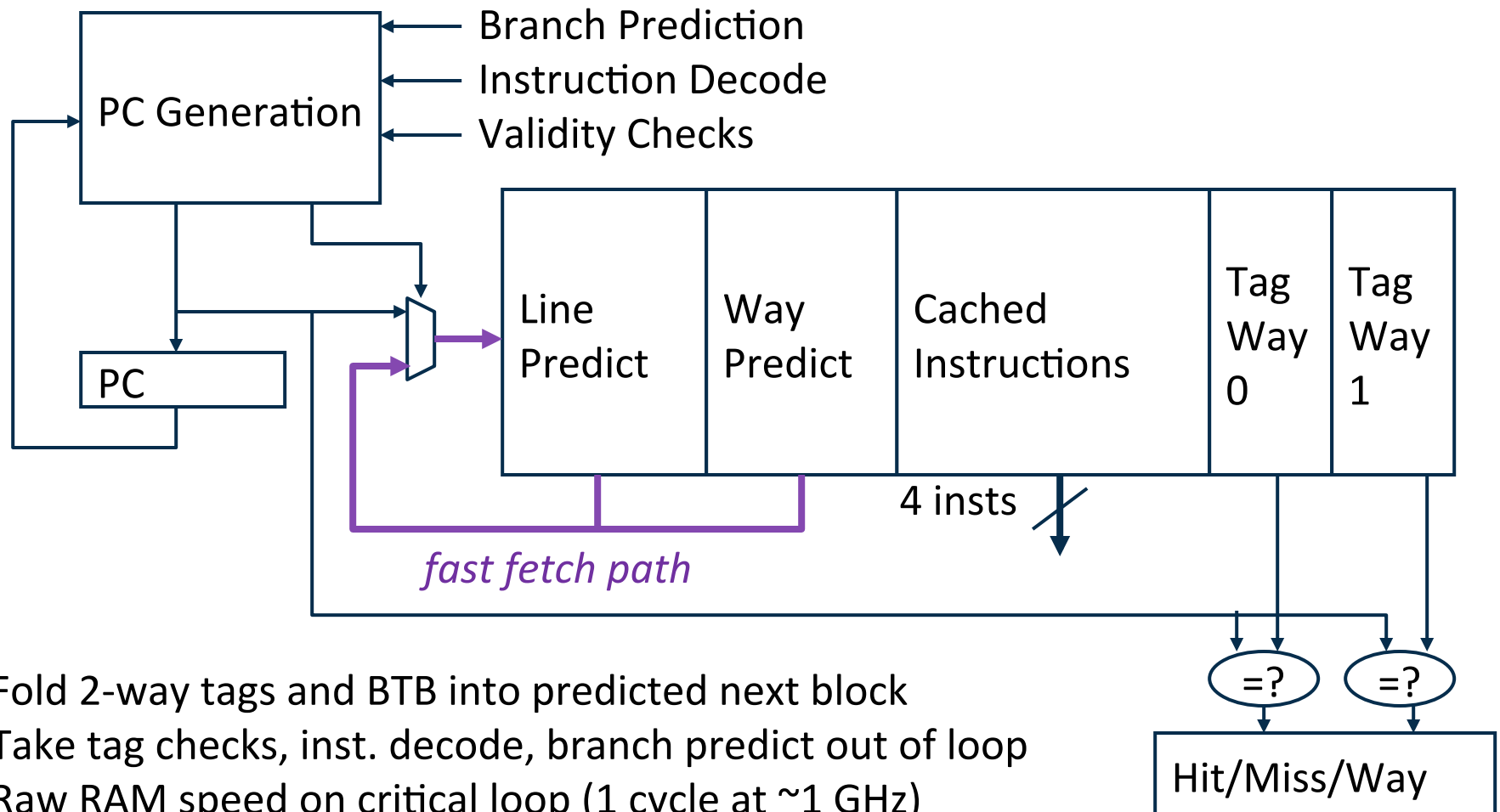
# Rename Table Recovery

- Have to quickly recover rename table on branch mispredicts
- MIPS R10K only has four snapshots for each of four outstanding speculative branches
- Alpha 21264 has 80 snapshots, one per ROB instruction

# Improving Instruction Fetch

- Performance of speculative out-of-order machines often limited by instruction fetch bandwidth
  - speculative execution can fetch 2-3x more instructions than are committed
  - mispredict penalties dominated by time to refill instruction window
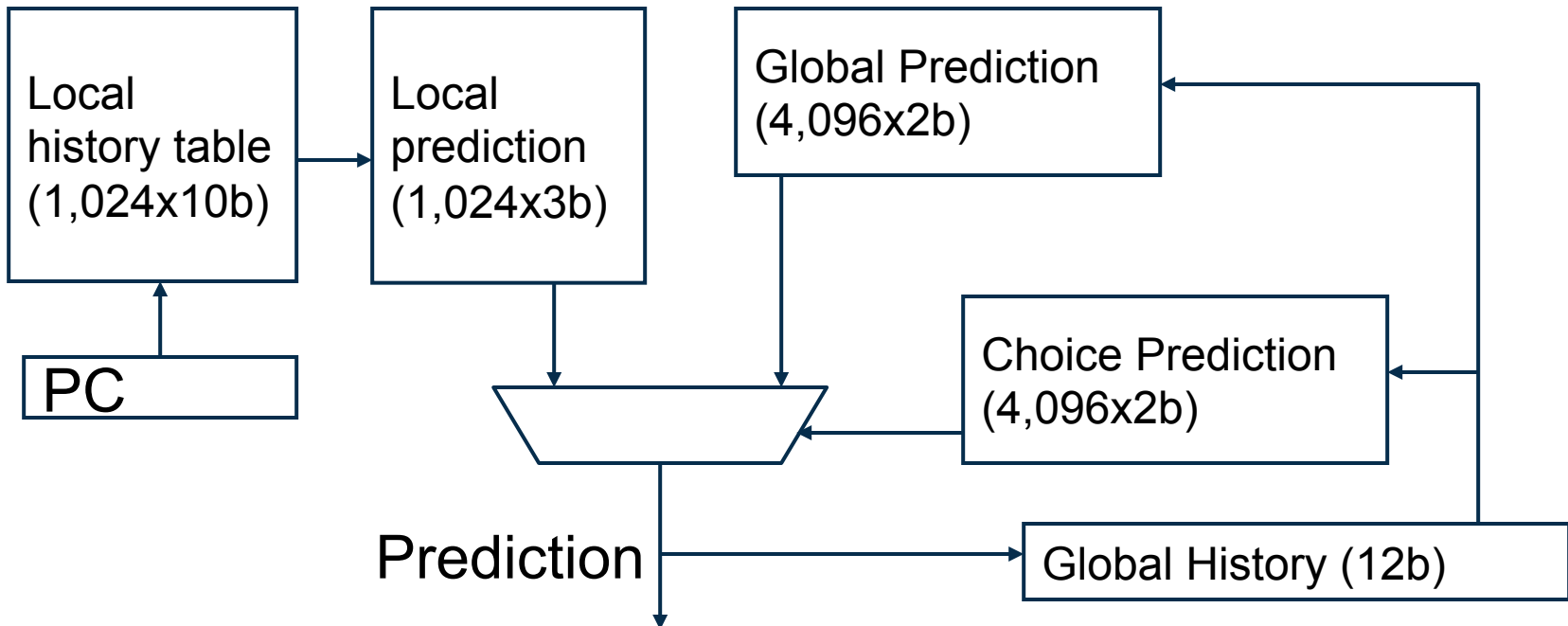  - taken branches are particularly troublesome

# Increasing Taken Branch Bandwidth
# (Alpha 21264 I-Cache)



- Fold 2-way tags and BTB into predicted next block
- Take tag checks, inst. decode, branch predict out of loop
- Raw RAM speed on critical loop (1 cycle at ~1 GHz)
- 2-bit hysteresis counter per block prevents overtraining

　　　　　© Krste Asanovic, 2015

# Tournament Branch Predictor
# (Alpha 21264)

- Choice predictor learns whether best to use local or global branch history in predicting next branch
- Global history is speculatively updated but restored on mispredict
- Claim 90-100% success on range of applications



```
Local                Local              Global Prediction
history table        prediction         (4,096x2b)
(1,024x10b)          (1,024x3b)

                                         Choice Prediction
PC                                       (4,096x2b)

           Prediction                    Global History (12b)
```
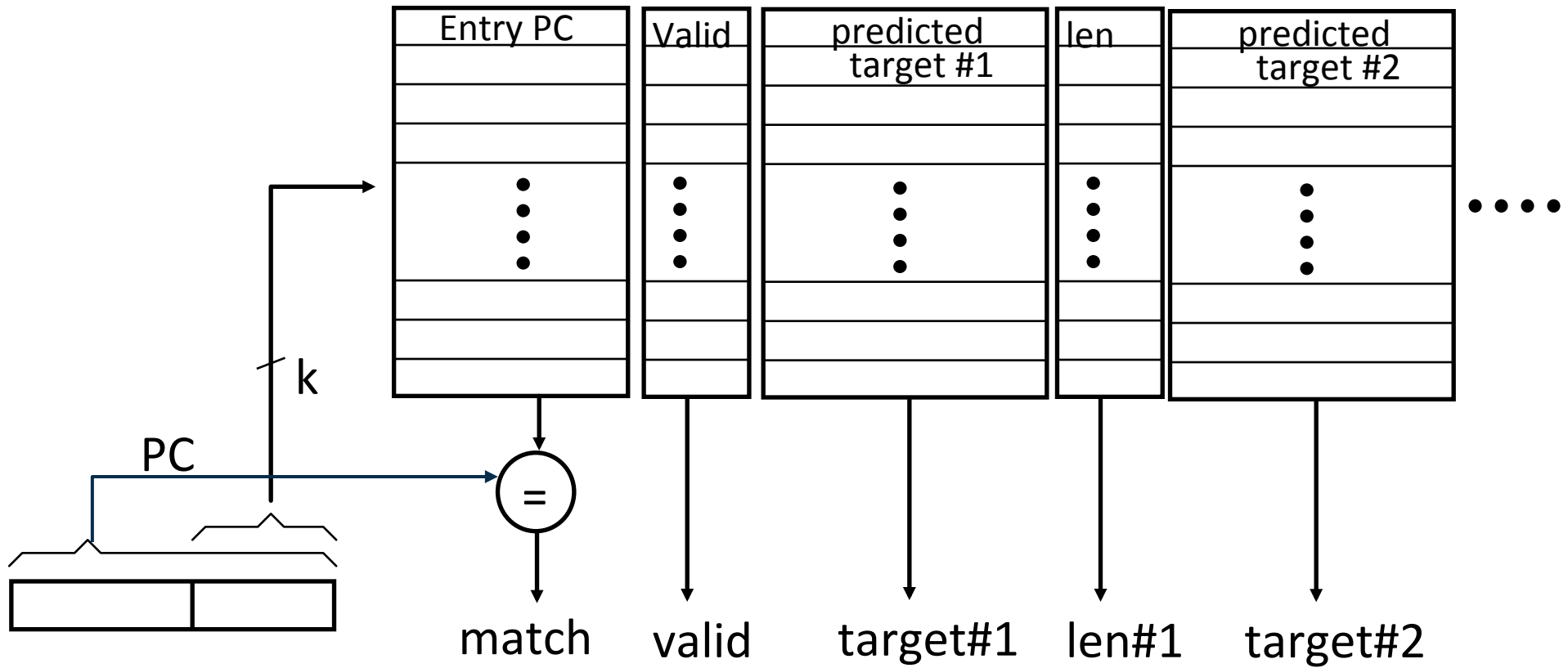
# Taken Branch Limit

- Integer codes have a taken branch every 6-9 instructions
- To avoid fetch bottleneck, must execute multiple taken branches per cycle when increasing performance
- This implies:
  - predicting multiple branches per cycle
  - fetching multiple non-contiguous blocks per cycle

# Branch Address Cache
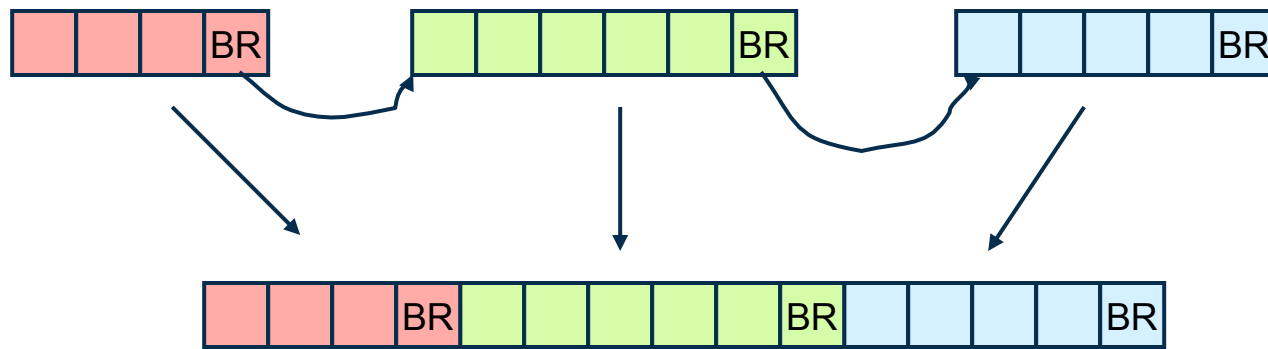# (Yeh, Marr, Patt)



Extend BTB to return multiple branch predictions per cycle

# Fetching Multiple Basic Blocks

- Requires either
  - multiported cache: expensive
  - interleaving: bank conflicts will occur

- Merging multiple blocks to feed to decoders adds latency increasing mispredict penalty and reducing branch throughput

© Krste Asanovic, 2015

# Trace Cache

- Key Idea: Pack multiple non-contiguous basic blocks into one contiguous trace cache line
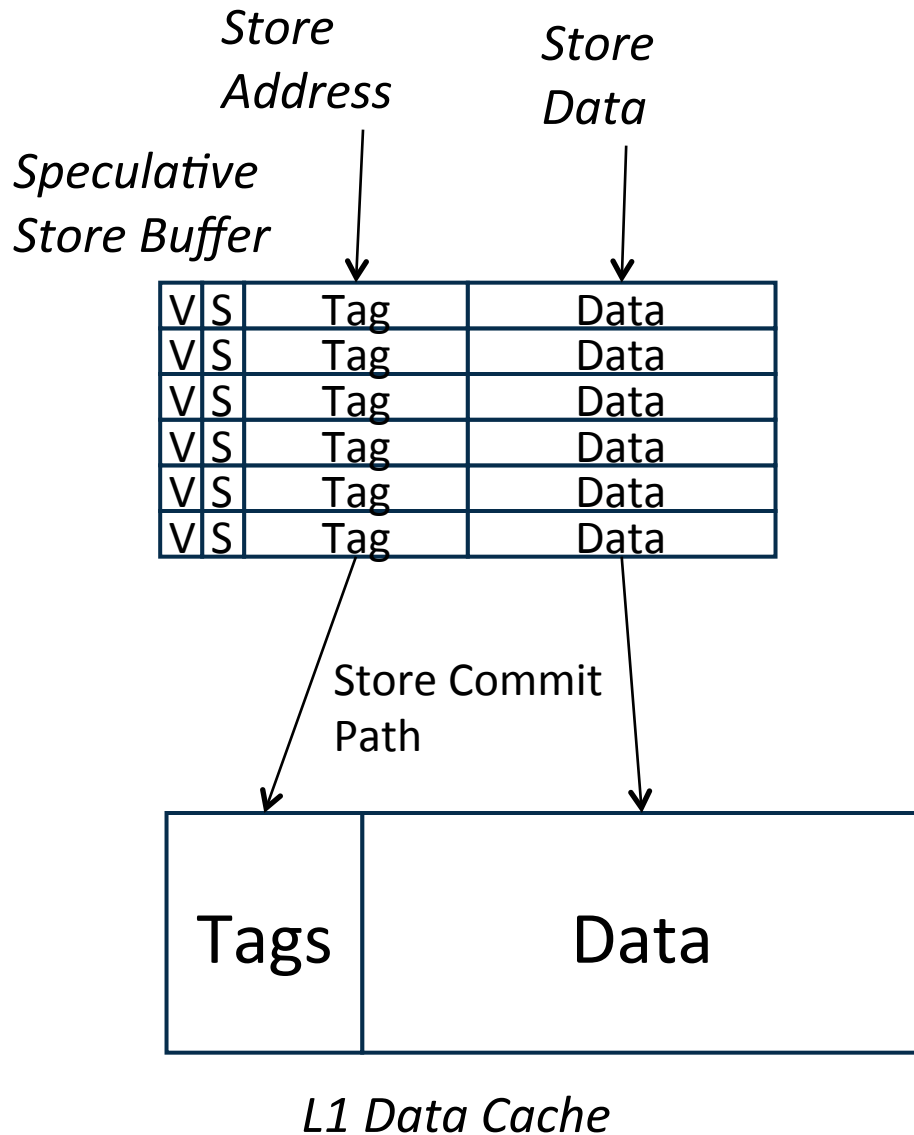


- Single fetch brings in multiple basic blocks

- Trace cache indexed by start address *and* next *n* branch predictions

- Used in Intel Pentium-4 processor to hold decoded uops
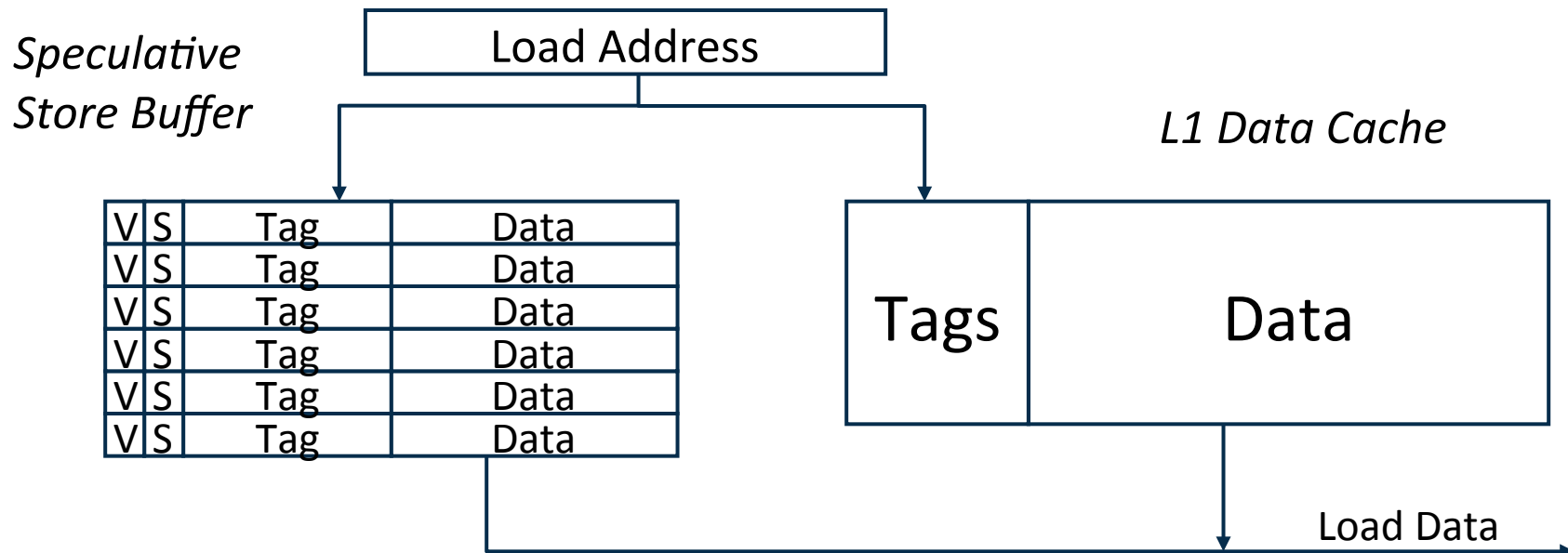
# Load-Store Queue Design

- After control hazards, data hazards through memory are probably next most important bottleneck to superscalar performance

- Modern superscalars use very sophisticated load-store reordering techniques to reduce effective memory latency by allowing loads to be speculatively issued

# Speculative Store Buffer

*Store Address*

*Store Data*

*Speculative Store Buffer*

| V | S | Tag | Data |
|---|---|-----|------|
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |

Store Commit Path

| Tags | Data |
|------|------|

*L1 Data Cache*

- Just like register updates, stores should not modify the memory until after the instruction is committed. A speculative store buffer is a structure introduced to hold speculative store data.
- During decode, store buffer slot allocated in program order
- Stores split into "store address" and "store data" micro-operations
- "Store address" execution writes tag
- "Store data" execution writes data
- Store commits when oldest instruction and both address and data available:
  - clear speculative bit and eventually move data to cache
- On store abort:
  - clear valid bit

# Load bypass from speculative store buffer

*Speculative Store Buffer*

Load Address

*L1 Data Cache*

| V | S | Tag | Data |
|---|---|-----|------|
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |
| V | S | Tag | Data |

| Tags | Data |
|------|------|

Load Data

- If data in both store buffer and cache, which should we use?
  Speculative store buffer
- If same address in store buffer twice, which should we use?
  Youngest store older than load

# Memory Dependencies

```
sd x1, (x2)
ld x3, (x4)
```

- When can we execute the load?

# In-Order Memory Queue

- Execute all loads and stores in program order

- => Load and store cannot leave ROB for execution until all previous loads and stores have completed execution

- Can still execute loads and stores speculatively, and out-of-order with respect to other instructions

- Need a structure to handle memory ordering...

# Conservative O-o-O Load Execution

```
sd x1, (x2)
ld x3, (x4)
```

- Can execute load before store, if addresses known and **x4** != **x2**
- Each load address compared with addresses of all previous uncommitted stores
  - can use partial conservative check i.e., bottom 12 bits of address, to save hardware
- Don't execute load if any previous store address not known
- (MIPS R10K, 16-entry address queue)

# Address Speculation

```
sd x1, (x2)
ld x3, (x4)
```

- Guess that **x4** != **x2**
- Execute load before store address known
- Need to hold all completed but uncommitted load/ store addresses in program order
- If subsequently find **x4==x2**, squash load and all following instructions


- => Large penalty for inaccurate address speculation

# Memory Dependence Prediction
# (Alpha 21264)

```
sd x1, (x2)
ld x3, (x4)
```

- Guess that **x4** != **x2** and execute load before store

- If later find **x4**==**x2**, squash load and all following instructions, but mark load instruction as store-wait

- Subsequent executions of the same load instruction will wait for all previous stores to complete

- Periodically clear store-wait bits

# Acknowledgements

- This course is partly inspired by previous MIT 6.823 and Berkeley CS252 computer architecture courses created by my collaborators and colleagues:
  - Arvind (MIT)
  - Joel Emer (Intel/MIT)
  - James Hoe (CMU)
  - John Kubiatowicz (UCB)
  - David Patterson (UCB)